

SPECIFIC EXPRESSION JUDGING DEVICE AND SPECIFIC EXPRESSION JUDGING METHOD, AND RECORDING MEDIUM WITH SPECIFIC EXPRESSION JUDGING PROGRAM RECORDED THEREON

Publication number: JP2002082943 (A)

Publication date: 2002-03-22

Inventor(s): FUKUSHIMA SHUNICHI +

Applicant(s): NIPPON ELECTRIC CO +

Classification:

- international: G06F17/22; G06F17/21; G06F17/30; G06F17/22; G06F17/21; G06F17/30; (IPC1-7): G06F17/22; G06F17/21

- European: G06F17/30G4; G06F17/30T

Application number: JP20000278691 20000908

Priority number(s): JP20000278691 20000908

Also published as:

JP4200645 (B2)

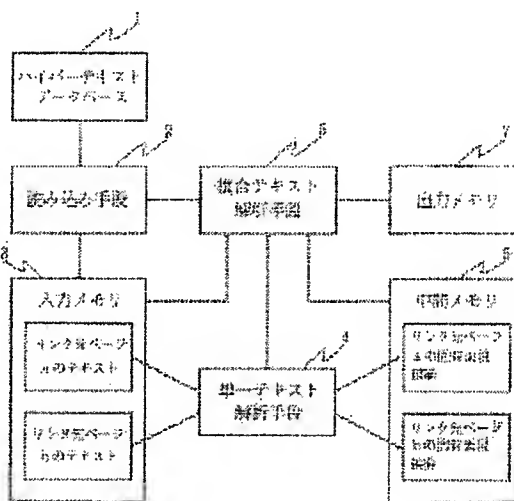
US2002031269 (A1)

US6975766 (B2)

Abstract of JP 2002082943 (A)

PROBLEM TO BE SOLVED: To provide a specific expression judging device capable of highly precisely judging a specific expression(a place name, a persons' name, an organization name or the like) appearing in the text of each node page constituting a hypertext data base such as a WWW.

SOLUTION: A reading means 2 reads a text from a hypertext data base 1. A single text analyzing means 4 detects a specific expression candidate appearing in each text read by the reading means 2 by the in-text analysis processing. A composite text analyzing means 6 calculates the certainty of a specific expression candidate detected by the single text analyzing means 4 according to the analysis processing for referring to the text at the origin of link or the destination of link corresponding to the text in which the specific expression candidate appears.



(19)日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11)特許出願公開番号
特開2002-82943
(P2002-82943A)

(43)公開日 平成14年3月22日(2002.3.22)

(51)Int.Cl. ⁷	識別記号	F I	テーマコード*(参考)
G 0 6 F 17/22	5 2 2	G 0 6 F 17/22	5 2 2 L 5 B 0 0 9
17/21	5 5 0	17/21	5 5 0 Z

審査請求 未請求 請求項の数21 O L (全 15 頁)

(21)出願番号 特願2000-278691(P2000-278691)

(22)出願日 平成12年9月8日(2000.9.8)

(71)出願人 000004237

日本電気株式会社

東京都港区芝五丁目7番1号

(72)発明者 福島 俊一

東京都港区芝五丁目7番1号 日本電気株式会社内

(74)代理人 100084250

弁理士 丸山 隆夫

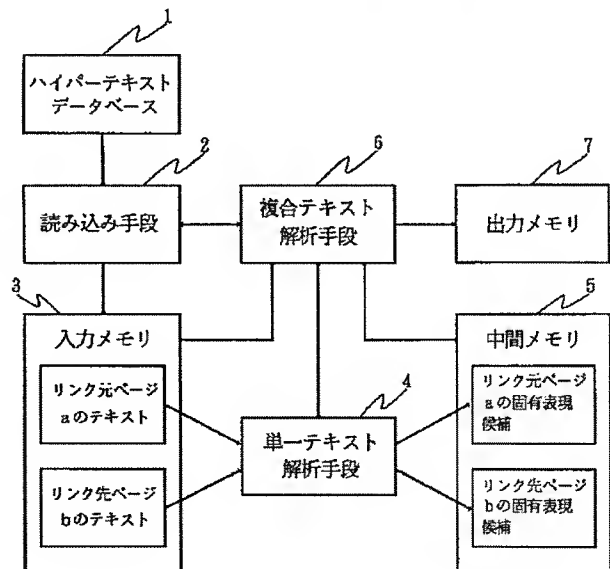
Fターム(参考) 5B009 MB07 MB21

(54)【発明の名称】 固有表現判別装置、固有表現判別方法、および固有表現判別プログラムを記録した記録媒体

(57)【要約】

【課題】 WWWのようなハイパーテキストデータベースを構成する各ノードページのテキスト中に出現する固有表現(地名・人名・組織名など)を高精度に判別可能な固有表現判別装置を提供する。

【解決手段】 読み込み手段2は、ハイパーテキストデータベース1からテキストを読み込む。単一テキスト解析手段4は、読み込み手段2によって読み込まれた各テキストから、そのテキスト内の解析処理によって、そのテキスト内に出現する固有表現候補を検出する。複合テキスト解析手段6は、単一テキスト解析手段4によって検出された固有表現候補の確からしさを、その固有表現候補の出現したテキストに対するリンク元あるいはリンク先のテキストを参照した解析処理によって計算する。



【特許請求の範囲】

【請求項 1】 テキスト中に出現する地名、人名、組織名などの固有表現を検出する固有表現判別装置において、

ハイパーテキストデータベースからテキストを読み込む読み込み手段と、

前記読み込み手段によって読み込まれた各テキストから、該テキスト内の解析処理によって、該テキスト内に出現する固有表現候補を検出する単一テキスト解析手段と、

前記単一テキスト解析手段によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストを参照した解析処理により計算する複合テキスト解析処理と、

を備えたことを特徴とする固有表現判別装置。

【請求項 2】 前記複合テキスト解析手段は、前記単一テキスト解析手段によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストにおける、該固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする請求項 1 記載の固有表現判別装置。

【請求項 3】 前記複合テキスト解析手段は、前記単一テキスト解析手段によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キーワード列と、該固有表現候補との共起関係を基に計算することを特徴とする請求項 1 記載の固有表現判別装置。

【請求項 4】 前記複合テキスト解析手段は、前記単一テキスト解析手段によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キーワード列の前後一定範囲のテキストにおける、該固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする請求項 1 記載の固有表現判別装置。

【請求項 5】 前記複合テキスト解析手段は、前記固有表現候補と共起する単語を前記リンク元、リンク先の両方若しくは一方のテキストから検出することができなかった場合には、該共起する単語を検出することができなかったテキストのリンク先若しくはリンク元のテキストを参照して、前記固有表現候補と共起する単語を検出することを特徴とする請求項 2 から 4 の何れか一項に記載の固有表現判別装置。

【請求項 6】 前記複合テキスト解析手段は、前記固有表現候補の出現したテキストからリンク元あるいはリンク先のテキストをたどり、該テキストから所定の階層までを参照範囲として前記固有表現候補と共起する単語を検出することを特徴とする請求項 2 から 4 の何れか一項に記載の固有表現判別装置。

【請求項 7】 前記複合テキスト解析手段は、前記固有表現候補と共起する単語の出現位置に応じて割り当てられた重みに基づき前記固有表現候補の確からしさを計算することを特徴とする請求項 2 から 4 の何れか一項に記載の固有表現判別装置。

【請求項 8】 ハイパーテキストデータベースからテキストを読み込み、読み込まれた各テキストから、該テキスト内の解析処理によって、該テキスト内に出現する固有表現候補を検出し、

前記固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストを参照した解析処理によって計算することを特徴とする固有表現判別方法。

【請求項 9】 前記固有表現候補の確からしさを計算する際に、

該固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストにおける、該固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする請求項 8 記載の固有表現判別方法。

【請求項 10】 前記固有表現候補の確からしさを計算する際に、

該固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キーワード列と、該固有表現候補との共起関係を基に計算することを特徴とする請求項 8 記載の固有表現判別方法。

【請求項 11】 前記固有表現候補の確からしさを計算する際に、

該固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キーワード列の前後一定範囲のテキストにおける、該固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする請求項 8 記載の固有表現判別方法。

【請求項 12】 前記固有表現候補の確からしさを計算する際に、

前記固有表現候補と共起する単語を前記リンク元、リンク先の両方若しくは一方のテキストから検出することができなかった場合には、該共起する単語を検出することができなかったテキストのリンク先若しくはリンク元のテキストを参照して、前記固有表現候補と共起する単語を検出することを特徴とする請求項 9 から 11 の何れか一項に記載の固有表現判別方法。

【請求項 13】 前記固有表現候補の確からしさを計算する際に、

前記固有表現候補の出現したテキストからリンク元あるいはリンク先のテキストをたどり、該テキストから所定の階層までを参照範囲として前記固有表現候補と共起する単語を検出することを特徴とする請求項 9 から 11 の何れか一項に記載の固有表現判別方法。

【請求項 14】 前記固有表現候補の確からしさを計算

する際に、

前記固有表現候補と共起する単語の出現位置に応じて割り当てられた重みに基づき前記固有表現候補の確からしさを計算することを特徴とする請求項 9 から 11 の何れか一項に記載の固有表現判別方法。

【請求項 15】 ハイパーテキストデータベースからテキストを読み込む読み込み処理と、

前記読み込み処理によって読み込まれた各テキストから、該テキスト内の解析処理によって、該テキスト内に出現する固有表現候補を検出する単一テキスト解析処理と、

前記単一テキスト解析手段によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストを参照した解析処理によって計算する複合テキスト解析処理と、

を実行するためのプログラムを記録したことを特徴とする固有表現判別プログラムを記録した記録媒体。

【請求項 16】 前記複合テキスト解析処理は、前記単一テキスト解析処理によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストにおける、該固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする請求項 15 記載の固有表現判別プログラムを記録した記録媒体。

【請求項 17】 前記複合テキスト解析処理は、前記単一テキスト解析処理によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キーワードと、該固有表現候補との共起関係を基に計算することを特徴とする請求項 15 記載の固有表現判別プログラムを記録した記録媒体。

【請求項 18】 前記複合テキスト解析処理は、前記単一テキスト解析処理によって検出された固有表現候補の確からしさを、該固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キーワードの前後一定範囲のテキストにおける、該固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする請求項 15 記載の固有表現判別プログラムを記録した記録媒体。

【請求項 19】 前記複合テキスト解析処理は、前記固有表現候補と共起する単語を前記リンク元、リンク先の両方若しくは一方のテキストから検出することができなかった場合には、該共起する単語を検出することができなかったテキストのリンク先若しくはリンク元のテキストを参照して、前記固有表現候補と共起する単語を検出することを特徴とする請求項 16 から 18 の何れか一項に記載の固有表現判別プログラムを記録した記録媒体。

【請求項 20】 前記複合テキスト解析処理は、

前記固有表現候補の出現したテキストからリンク元あるいはリンク先のテキストをたどり、該テキストから所定の階層までを参照範囲として前記固有表現候補と共起する単語を検出することを特徴とする請求項 16 から 18 の何れか一項に記載の固有表現判別プログラムを記録した記録媒体。

【請求項 21】 前記複合テキスト解析処理は、前記固有表現候補と共起する単語の出現位置に応じて割り当てられた重みに基づき前記固有表現候補の確からしさを計算することを特徴とする請求項 16 から 18 の何れか一項に記載の固有表現判別プログラムを記録した記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、テキスト中出现する地名・人名・組織名などの固有表現を検出する固有表現判別装置、固有表現判別方法、および、固有表現判別プログラムを記録した記録媒体に関する。ここで、固有表現は、Named Entity に対応する日本語であり、地名・人名・組織名などを指す。「言語処理学会第 5 回年次大会」論文集（1999 年 3 月）の pp.128~131 に掲載された論文「固有表現の定義の困難さーIREXにおけるNE定義の事例からー」（著者：関根聡・江里口善生）、あるいは、1999 年 9 月に開催された「IREXワークショップ」などにおいて当該分野の用語として定義されている。

【0002】

【従来の技術】固有表現を検出するための最も基本的な方式は、固有表現の辞書を用意し、テキストと辞書とを照合することで、テキスト中出现した固有表現を検出するものである。例えば、辞書のなかに「横浜市」（地名）、「横浜ベイスターズ」（組織名）のように登録しておき、テキスト中に「横浜市」が出現すれば、それを地名として検出し、「横浜ベイスターズ」が出現すれば、それを組織名として検出する。

【0003】しかし、単純に辞書と照合するだけでは、固有表現を判別できないことがある。例えば、テキスト中に「千葉」という表記が出現した場合、これは人名かもしれないし、地名かもしれないという、複数通りの解釈（曖昧性）が生ずる。同様に、テキスト中に「谷」という表記が出現した場合、これは人名かもしれないし、一般名詞かもしれないという曖昧性を持つ。さらには、テキスト中の「中央区」という表記が地名として検出できたとしても、この「中央区」が、「東京都中央区」なのか、「大阪府中央区」なのか、という解釈の曖昧性は残る。

【0004】このような固有表現の判別における曖昧性を解消するための手法として、従来、以下のような 2 通りの方法が考えられている。これらの手法はいずれも、「IREXワークショップ」（1999 年 9 月）の予稿

集に掲載された論文群、特に、「固有表現抽出システムの開発とIREX-NEにおける評価」（著者：竹元義美・福島俊一・山田洋志・奥村明俊・池田崇博）などに記載されている。

【0005】第一の手法は、固有表現の候補の前後あるいは同一テキスト内に出現する共起語を参照して、固有表現候補の曖昧性を解消する方法である。例えば、「千葉」という固有表現候補の直後に「選手」という共起語が出現すれば、この「千葉」は人名と判定できる。あるいは、「中央区」という固有表現候補について、同じ

10

テキスト中に「東京都」という共起語が出現していれば、この「中央区」は「東京都中央区」を意味する可能性が高いと解釈できる。

【0006】第二の手法は、固有表現の候補の表記を包含するような表記が、同一テキスト内に出現しているかを調べて、固有表現候補の曖昧性を解消する方法である。例えば、「横浜」という表記には地名と組織名の曖昧性があるが、同一テキスト内に「横浜ベイスターズ」という表記が出現しているならば、「横浜」は「横浜ベイスターズ」の省略表記、すなわち、組織名である可能

20

性が高いと判断できる。

【0007】本発明と技術分野が類似する従来例1として、特開平6-52221号公報の“固有名詞の自動抽出方式”がある。

【0008】本従来例は、オンライン・データベースやエキスパート・システム、機械翻訳システム等の自然言語インターフェースにおいて、オンライン・テキストをアクセスするデータベース・アクセス手段と、前もって作成してある固有名詞パターンを参照し、データベース・アクセス手段がアクセスして得たテキストから固有名

30

詞候補を抽出する固有名詞抽出手段と、固有名詞抽出手段が抽出した固有名詞候補が既に辞書に登録してあるか否かを判定し、未登録の固有名詞候補を抽出する固有名詞判定手段と、固有名詞判定手段が抽出した未登録の固有名詞候補を辞書に登録する固有名詞登録手段とを有することを特徴としている。

【0009】

【発明が解決しようとする課題】しかしながら、上述した固有表現判別方式、並びに従来例1の固有名詞の自動抽出方式は、1つのテキスト内の解析処理によるものである。このような従来技術では、WWW (World Wide Web) のようなハイパーテキストデータベースを構成する各ノードページのテキストを対象とした場合、そのテキスト内の情報のみを用いた解析処理では、十分な判別精度が得られない可能性があるという問題がある。

40

【0010】本発明は上記事情に鑑みてなされたものであり、WWWのようなハイパーテキストデータベースを構成する各ノードページのテキスト中に出現する固有表現（地名・人名・組織名など）を高精度に判別可能な固有表現判別装置、固有表現判別方法、固有表現判別プロ

50

グラムを記録した記録媒体を提供することを目的とする。

【0011】

【課題を解決するための手段】係る目的を達成するために請求項1記載の発明は、テキスト中に出現する地名、人名、組織名などの固有表現を検出する固有表現判別装置において、ハイパーテキストデータベースからテキストを読み込む読み込み手段と、読み込み手段によって読み込まれた各テキストから、テキスト内の解析処理によって、テキスト内に出現する固有表現候補を検出する単一テキスト解析手段と、単一テキスト解析手段によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストを参照した解析処理により計算する複合テキスト解析処理と、を備えたことを特徴とする。

【0012】請求項2記載の発明は、請求項1記載の発明において、複合テキスト解析手段は、単一テキスト解析手段によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストにおける、固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする。

【0013】請求項3記載の発明は、請求項1記載の発明において、複合テキスト解析手段は、単一テキスト解析手段によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列と、固有表現候補との共起関係を基に計算することを特徴とする。

【0014】請求項4記載の発明は、請求項1記載の発明において、複合テキスト解析手段は、単一テキスト解析手段によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列の前後一定範囲のテキストにおける、固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする。

【0015】請求項5記載の発明は、請求項2から4の何れか一項に記載の発明において、複合テキスト解析手段は、固有表現候補と共起する単語をリンク元、リンク先の両方若しくは一方のテキストから検出することができなかった場合には、共起する単語を検出することができなかったテキストのリンク先若しくはリンク元のテキストを参照して、固有表現候補と共起する単語を検出することを特徴とする。

【0016】請求項6記載の発明は、請求項2から4の何れか一項に記載の発明において、複合テキスト解析手段は、固有表現候補の出現したテキストからリンク元あるいはリンク先のテキストをたどり、テキストから所定の階層までを参照範囲として固有表現候補と共起する単語を検出することを特徴とする。

【0017】請求項7記載の発明は、請求項2から4の何れか一項に記載の発明において、複合テキスト解析手段は、固有表現候補と共起する単語の出現位置に応じて割り当てられた重みに基づき固有表現候補の確からしさを計算することを特徴とする。

【0018】請求項8記載の発明は、ハイパーテキストデータベースからテキストを読み込み、読み込まれた各テキストから、テキスト内の解析処理によって、テキスト内に出現する固有表現候補を検出し、固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストを参照した解析処理によって計算することを特徴とする。

【0019】請求項9記載の発明は、請求項8記載の発明において、固有表現候補の確からしさを計算する際に、固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストにおける、固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする。

【0020】請求項10記載の発明は、請求項8記載の発明において、固有表現候補の確からしさを計算する際に、固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列と、固有表現候補との共起関係を基に計算することを特徴とする。

【0021】請求項11記載の発明は、請求項8記載の発明において、固有表現候補の確からしさを計算する際に、固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列の前後一定範囲のテキストにおける、固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする。

【0022】請求項12記載の発明は、請求項9から11の何れか一項に記載の発明において、固有表現候補の確からしさを計算する際に、固有表現候補と共起する単語をリンク元、リンク先の両方若しくは一方のテキストから検出することができなかった場合には、共起する単語を検出することができなかったテキストのリンク先若しくはリンク元のテキストを参照して、固有表現候補と共起する単語を検出することを特徴とする。

【0023】請求項13記載の発明は、請求項9から11の何れか一項に記載の発明において、固有表現候補の確からしさを計算する際に、固有表現候補の出現したテキストからリンク元あるいはリンク先のテキストをたどり、テキストから所定の階層までを参照範囲として固有表現候補と共起する単語を検出することを特徴とする。

【0024】請求項14記載の発明は、請求項9から11の何れか一項に記載の発明において、固有表現候補の確からしさを計算する際に、固有表現候補と共起する単語の出現位置に応じて割り当てられた重みに基づき固有表現候補の確からしさを計算することを特徴とする。

【0025】請求項15記載の発明は、ハイパーテキストデータベースからテキストを読み込む読み込み処理

と、読み込み処理によって読み込まれた各テキストから、テキスト内の解析処理によって、テキスト内に出現する固有表現候補を検出する単一テキスト解析処理と、単一テキスト解析手段によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストを参照した解析処理によって計算する複合テキスト解析処理と、を実行するためのプログラムを記録したことを特徴とする。

【0026】請求項16記載の発明は、請求項15記載の発明において、複合テキスト解析処理は、単一テキスト解析処理によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元、リンク先の両方若しくは一方のテキストにおける、固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする。

【0027】請求項17記載の発明は、請求項15記載の発明において、複合テキスト解析処理は、単一テキスト解析処理によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列と、固有表現候補との共起関係を基に計算することを特徴とする。

【0028】請求項18記載の発明は、請求項15記載の発明において、複合テキスト解析処理は、単一テキスト解析処理によって検出された固有表現候補の確からしさを、固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列の前後一定範囲のテキストにおける、固有表現候補と共起する単語の出現頻度情報を基に計算することを特徴とする。

【0029】請求項19記載の発明は、請求項16から18の何れか一項に記載の発明において、複合テキスト解析処理は、固有表現候補と共起する単語をリンク元、リンク先の両方若しくは一方のテキストから検出することができなかった場合には、共起する単語を検出することができなかったテキストのリンク先若しくはリンク元のテキストを参照して、固有表現候補と共起する単語を検出することを特徴とする。

【0030】請求項20記載の発明は、請求項16から18の何れか一項に記載の発明において、複合テキスト解析処理は、固有表現候補の出現したテキストからリンク元あるいはリンク先のテキストをたどり、テキストから所定の階層までを参照範囲として固有表現候補と共起する単語を検出することを特徴とする。

【0031】請求項21記載の発明は、請求項16から18の何れか一項に記載の発明において、複合テキスト解析処理は、固有表現候補と共起する単語の出現位置に応じて割り当てられた重みに基づき固有表現候補の確からしさを計算することを特徴とする。

【0032】

【発明の実施の形態】次に、添付図面を参照しながら本

発明に係る実施の形態を詳細に説明する。図1～図16を参照すると本発明に係る実施の形態が示されている。

【0033】本発明に係る第1の実施形態は、図1に示されるように、読み込み手段2、入力メモリ3、単一テキスト解析手段4、中間メモリ5、複合テキスト解析手段6、出力メモリ7を備え、ハイパーテキストデータベース1を構成するノードページのテキストに対して、そのテキスト中に出現する固有表現のリストを出力する。これらの各手段は、プログラム制御によって動作するコンピュータを用いて実現できる。入力メモリ3、中間メモリ5、出力メモリ7は、コンピュータの主記憶部を用いてもよいが、磁気ディスク装置や光磁気ディスク装置などの外部記憶装置を用いてもよい。

【0034】まず、本実施形態が処理対象とするハイパーテキストについて説明する。

【0035】図2は、固有表現判別装置の入力となるハイパーテキストデータベース1の一部を示す例である。ハイパーテキストは、ページ（あるいはノードと呼ばれる）をテキストの1単位として、それらの間にリンクが設けられた形式をしている。図2における10と11は各々、ハイパーテキストデータベースを構成する1ページであり、それらの間をつなぐ矢印13はリンクを表している。リンク13に着目するならば、ページ10はリンク元ページ（リンク元テキスト）、ページ11はリンク先ページ（リンク先テキスト）となる。このようなハイパーテキストは、独自のデータ構造をもつものもあるが、最近ではSGML（Standard Generalized Markup Language）、HTML（Hyper Text Markup Language）、XML（Extensible Markup Language）などのマークアップ言語による記述が普及している。特に、インターネット上にはWWWと呼ばれる大規模ハイパーテキストデータベースが存在し、そのなかではHTMLによる記述がスタンダードになっている。

【0036】図3は、図2のページ10をHTMLで記述した一例である。図3において、<>で囲まれた部分はマークアップタグであり、HTMLではAタグで囲まれた文字列がリンク元を表す。すなわち、図2における文字列「ベイスターズファン」は、リンク元キー文字列になる。さらに、AタグのなかでHREFの直後に書かれたHTMLテキスト名が、そのリンク先のページ（テキスト）を意味する。すなわち、ページ10のリンク元キー文字列「ベイスターズファン」から「Baystars Fan」という名前のHTMLテキスト（ページ11がこれに相当する）へジャンプできることを意味している。なお、ここではHTMLで記述されたハイパーテキストデータベースを例にあげて説明したが、本発明では、対象とするハイパーテキストの記述形式をHTMLに限定するものではない。SGMLやXMLで記述されたものでもよいし、独自のデータ構造を用いたハイパーテキストであつてもかまわない。

【0037】読み込み手段2は、ハイパーテキストデータベース1から入力メモリ3へページ（テキスト）を読み込む。この読み込み手段2は、ハイパーテキストデータベース1がどこに置かれているかに応じて、ネットワークを介して外部と通信するための機構、あるいは、外部記憶装置にアクセスするための機構なども含む。どのページを読み込むかについては、（a）その都度、読み込む対象ページを指定して個別に読み込む方法、（b）ハイパーテキストの一部分を読み込む対象として範囲指定する方法、（c）ハイパーテキストのリンクを自動的にたどりながら、すべてのページを読み込む方法、などが考えられる。（a）や（b）の場合、外部から対象ページあるいは対象範囲を指定するため、読み込み手段2は、キーボードやマウスなど外部からコンピュータに指示入力を与える装置も含む。このような読み込み手段2の実現方法は、既に公知である。例えば、HTML形式のハイパーテキストデータベースを対象とするならば、

（a）はNetscape NavigatorやMicrosoft Internet Explorerなど広く普及したWWWブラウザのもつ基本機能であり、（b）はオートパイロットツールと呼ばれるものの、（c）はWWWロボットやクローラと呼ばれるシステムとして実用化されている。

【0038】単一テキスト解析手段4は、入力メモリ3に読み込まれた各テキストから、そのテキスト内の解析処理によって、そのテキスト内に出現する固有表現候補を検出する。この単一テキスト解析手段4は、従来の固有表現判別装置に相当する。すなわち、本明細書の従来技術の項で述べたような公知の技術によって実現できる。単一テキスト解析手段4は、入力メモリ3からテキストを1件読み込み、そのテキストに対する解析処理を実行し、そのテキストにおける固有表現候補を中間メモリ5へ書き込む。中間メモリ5に書き込まれる情報は、対象テキスト中に検出された固有表現候補のリストであり、個々の固有表現候補に関しては、その固有表現候補がテキスト中のどこに出現したかという位置情報、および、その固有表現候補の種類（地名、人名、組織名などのいずれか）などを記録する。その際、固有表現候補の種類に曖昧性がある場合には、地名または人名、人名または一般名詞、というように複数通りの可能性を記録しておく。なお、図4には、単一テキスト解析手段4によって検出した固有表現候補の情報を記録する中間メモリ5の管理テーブルの構成が示されている。

【0039】複合テキスト解析手段6は、単一テキスト解析手段4によって検出された固有表現候補の確からしさを、その固有表現候補の出現したリンク元あるいはリンク先のテキストを参照した解析処理によって計算する。この解析処理の実現方法、すなわち、リンク元あるいはリンク先のテキストを参照することによる固有表現候補の確からしさの計算方法としては、例えば、以下のようないものが考えられる。（ア）固有表現候補の確か

しさを、その固有表現候補の出現したテキストに対するリンク元あるいはリンク先の複数テキストにおける、その固有表現候補と共起する単語の出現頻度情報をもとにして計算する方法、(イ) 固有表現候補の確からしさを、その固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列と、その固有表現候補との共起関係をもとにして計算する方法、

(ウ) 固有表現候補の確からしさを、その固有表現候補の出現したテキストに対するリンク元テキストにおけるリンク元キー文字列の前後一定範囲のテキストにおける、その固有表現候補と共起する単語の出現頻度情報をもとにして計算する方法、などである。このような計算を実行するために、複合テキスト解析手段 6 は、必要に応じて、入力メモリ 3 や中間メモリ 5 の内容を参照する。そして、複合テキスト解析手段 6 の結果は、出力メモリ 7 へ書き出す。また、図 1 における複合テキスト解析手段 6 は、装置全体を制御するような役割を持たされた構成になっており、読み込み手段 2 や単一テキスト解析手段 4 の動作制御も行う。ただし、そのような装置全体の動作制御の機能は、複合テキスト解析手段 6 と分離して構成するようにしてもよい。

【0040】上記構成からなる本実施形態は、WWW のようなハイパーテキストデータベースを構成するノードページのテキスト中に出現する固有表現を高精度に判別することを目的としている。

【0041】例えば、図 2 に示されたテキストには、「横浜」という表記が出現する。「横浜」という表記には地名と組織名の曖昧性があり、そのテキストを検索しただけでは、この曖昧性を解消することができない場合がある。また、図 5 では、テキスト 16 中に出現する地名「中央区」が「大阪市中央区」なのか、「東京都中央区」なのかといった曖昧性が生じている。

【0042】このような不具合を解決するために、本実施形態は、単一テキスト解析手段 4 によって検出したテキスト内に出現する固有表現候補の確からしさを、その固有表現候補の出現したテキストに対するリンク元、あるいはリンク先のテキストを参照した解析によって計算する複合テキスト解析手段 6 を設けたことを特徴としている。この複合テキスト解析手段 6 により、例えば、図 2 に示された例では、固有表現候補「横浜」が現れたテキストのリンク元テキスト 10 を参照することで、「横浜」は組織名である方が確からしいと判定することができる。また、図 5 に示された例では、リンク元テキスト 15 を参照すると、「大阪府」という共起語が出現していることから、テキスト 16 の「中央区」は「大阪市中央区」の解釈の方が確からしいと判定できる。

【0043】図 6 を参照しながら、本発明の特徴部分である複合テキスト解析手段 6 の詳細な構成及び動作を説明する。図 6 に示されるように複合テキスト解析手段 6 は、入力メモリ読み込み手段 30、テキストバッファ 3

1、固有表現候補バッファ 32、固有表現辞書部 33、共起語情報読み出し部 34、共起語情報バッファ 35、共起語検出部 36、検出結果記憶部 37、尤度計算部 38 を有して構成される。

【0044】入力メモリ読み込み手段 30 は、入力メモリ 3 に記憶されたテキストを読み込む。この入力メモリ読み込み手段 30 の読み込み対象となるテキストの範囲は、上述した (ア)、(イ)、(ウ) のそれぞれの方法において異なるが、これらについては後に詳述する。テキストバッファ 31 は、入力メモリ読み込み部 30 により読み込まれたテキストを一時的に記憶する。

【0045】固有表現候補バッファ 32 は、単一テキスト解析手段 4 の解析結果である固有表現候補を中間メモリ 5 から読み出す。

【0046】固有表現辞書 33 には、固有表現候補を特定するための辞書が記憶されている。図 7 にこの固有表現辞書の構成を示す。図 7 に示されるように固有表現辞書では、固有表現の表記 40 に対して、地名、人名、組織名などのいずれかの固有表現として解釈されるかの種別 41 (一般名詞との解釈の曖昧性がある場合など固有表現以外の種別を入れることもある。) 及び各々の解釈時の共起語リスト 42 を格納している。共起語リスト 42 は、単に共起語のリストだけではなく、位置関係の条件 (固有表現の直後に接続するなど) も併せて格納しておいてもよい。

【0047】共起語情報読み出し部 34 は、固有表現候補バッファ 32 から固有表現候補、その位置情報、種類を読み出すと共に、この固有表現候補の、種別 41、及びその共起語リスト 42 を固有表現辞書 33 から読み出す。共起語情報読み出し部 34 により読み出された情報は、共起語情報バッファ 35 に一時的に記憶される。

【0048】共起語検出部 36 は、テキストバッファ 31 からテキストを読み出すと共に、共起語情報バッファ 35 から固有表現候補の共起語リストを読み出し、テキストの中から固有表現候補の共起語リストに挙げられた共起語を検出する。共起語検出部 36 にて検出された検出結果は、検出結果記憶部 37 に記憶される。

【0049】尤度計算部 38 は、検出結果記憶部 37 に記憶された共起語から固有表現候補の種別 (地名、人名、組織名) 毎の確からしさを判定し、判定結果を出力メモリ 7 に出力する。

【0050】ここで、複合テキスト解析手段 6 による (ア)、(イ)、(ウ) の解析手法について図 2 に示されたリンク元、及びリンク先のテキストを参照しながら説明すると共に、上記各解析手法の場合の複合テキスト解析手段 6 の処理動作を説明する。

【0051】図 2 では、10 がリンク元テキスト、11 がリンク先テキストである。リンク先テキスト 11 を対象として単一テキスト解析手段 4 が実行され、「横浜」が固有表現候補として検出されたものとする。図 7 に示

された固有表現辞書を用いたとするならば、図2のリンク先テキスト11における「横浜」には、地名と組織名という2通りの解釈（曖昧性）が生ずる。この曖昧性は、「横浜」に関する共起語が、テキスト11内に出現していないことから、単一テキスト解析手段4では解消できない。このとき、複合テキスト解析手段6では、リンク元テキスト10を参照することで、テキスト11における固有表現候補「横浜」の確からしさを計算する。複合テキスト解析手段6の実現方法として前述の

(ア)、(イ)、(ウ)では、次のように参照するテキスト範囲が異なる。

【0052】まず(ア)では、リンク元テキスト10の全体を参照する。そして、そのなか中出现する「横浜」の共起語として、「プロ野球」、「球団」、「ベ이스ターズ」などを見つけることで、固有表現候補「横浜」は組織名としての解釈の方が確からしいことを判定できる。

【0053】この解析方法の場合、複合テキスト解析手段6の入力メモリ読み込み部30は、入力メモリ3に記憶されたテキストの中からリンク元ページのテキスト全体を読み込み、テキストバッファ31に記憶する。

【0054】共起語検出部36は、テキストバッファ31からリンク元ページのテキスト全体を参照して、共起語情報バッファ35から読み出した共起語リストに挙げられた固有表現候補の共起語を検出する。

【0055】次に、(イ)の解析方法の場合、リンク元キー文字列のみを参照する。図2では、12の「ベ이스ターズファン」がリンク元キー文字列となる。この12のなかに、「横浜」の共起語である「ベ이스ターズ」が出現していることから、固有表現候補「横浜」は組織名としての解釈の方が確からしいことを判定できる。

【0056】この解析方法の場合、複合テキスト解析手段6の入力メモリ読み込み部30は、入力メモリ3に記憶されたリンク元ページのテキストの中から、リンク元キー文字列だけを読み込む。すなわち、リンク元ページの中から、固有表現候補の出現するテキストのテキスト名を、リンク先テキストに設定された文字列を読み込む。図2に示された例では、固有表現候補「横浜」の現れるテキスト名「Baystars Fan. html」が、アンカータグ(……)内に記載されたリンク元キー文字列を入力する。

【0057】共起語検出部36は、テキストバッファ31から読み込んだこのリンク元キー文字列に、共起語情報バッファ35から読み込んだ共起語リストに挙げられた固有表現候補の共起語が現れるか否かを検出する。

【0058】次に、(ウ)の解析方法の場合、図8における14のように、リンク元キー文字列の前後一定範囲のテキストを参照する。図8の14の範囲には、「球団」「ベ이스ターズ」などが「横浜」の共起語として出現しているので、固有表現候補「横浜」は組織名として

の解釈の方が確からしいことを判定できる。なお、リンク元キー文字列の前後一定範囲の決め方は、前後に一定の文字数、前後に一定の行数、リンク元キー文字列を含む1段落(～3段落)などのように、いろいろな方法が考えられる。

【0059】この解析方法の場合、複合テキスト解析手段6の入力メモリ読み込み部30は、入力メモリ3に記憶されたリンク元ページのテキストの中から、リンク元キー文字列、及びこのリンク元キー文字列の前後を一定の範囲を読み込む。

【0060】共起語検出部36は、テキストバッファ31から読み込んだこのリンク元キー文字列、及びこの前後一定の範囲のテキスト内に、共起語情報バッファ35から読み込んだ共起語リストに挙げられた固有表現候補の共起語が現れるか否かを検出する。

【0061】このようにして、本実施形態は、固有表現候補の出現したテキスト内だけではなく、リンク元のテキストも参照して固有表現候補の種別を特定することで、より高精度な固有表現の判別を行うことができる。

【0062】なお、複数テキストにおける共起語の出現頻度情報に着目する際の計算方法には、いろいろなバリエーションが考えられる。例えば、図9のテキスト17とテキスト19を考えて見ると、テキスト19における「中央区」の曖昧性に対して、リンク元テキストである17には「東京都」と「大阪府」の両方が出現していて、曖昧性を解消できない。そこで、(ア)の方法では、リンク元テキスト1件だけでなく複数件を参照する。さらには、リンク先テキストも参照することまで行うようにしている。図9のテキスト19に対するリンク元テキスト17と18、さらにリンク先テキスト20を参照すると、「東京都」(1回)、「大阪府」(3回)、「近畿地方」(1回)、「京都府」(1回)などが出現しており、最も多く出現している共起語である「大阪府」に着目することで、「中央区」は「大阪市中心区」という解釈の方が確からしいと判断できる。

【0063】また、上述した方法では、固有表現候補の曖昧性解消の際に、共起語のうちでリンク元・リンク先の複数テキストにおける出現頻度の総和が最も大きいものを優先した。それ以外にも、共起語のうちでリンク元・リンク先のなるべく多数件のテキストに出現するものを優先する方法も考えられる。これを図9の例で説明すれば、「大阪府」の出現するテキストは17・18・20の3件、「近畿地方」の出現するテキストは18のみで1件、「京都府」の出現するテキストも18のみで1件とカウントし、「大阪府」が最も多数件のテキストに出現した共起語ということになり、これを曖昧性解消の手がかりに用いるという方法である。

【0064】さらに、単語(共起語)の出現回数や出現テキスト件数の単純カウントではなく、リンク元テキストか、リンク先テキストかによって、重みを変えてカウ

ントするという方法も考えられる。例えば、単語の出現回数をリンク元テキストについては2点、リンク先テキストについては1点というように異なる重みを与えると、図9の例に関して、「東京都」は2点、「大阪府」は5点、「近畿地方」は2点、「京都府」は2点となる。また、共起語の出現位置によって重みを変える方法も考えられる。例えば、リンク元キー文字列に出現した場合は4点、リンク元キー文字列の前後一定範囲のテキストに出現した場合は3点、リンク元テキスト内に出現した場合は2点、リンク先テキスト内に出現した場合は1点というような重みの付け方である。

【0065】また、以上で述べた複合テキスト解析手段6の実現方法・処理例では、固有表現候補の出現したテキストからリンクを1階層分たどった範囲で、曖昧性解消の手がかりとなる共起語を探した。しかし、1階層よりも広い範囲から共起語を探すようにしてもよい。図10の例では、テキスト23に出現した「横浜」が固有表現候補であり、地名と組織名という曖昧性を持つ。このテキスト23のリンク元テキストの範囲（1階層分）では、「横浜」の曖昧性を解消する手がかりとなる共起語は出現していない。そこで、さらにもう1階層分、リンク元へ逆上ると、テキスト21を参照できる。テキスト21には「ベ이스ターズ」という共起語が存在するため、「横浜」は組織名としての解釈を優先することができる。このような複数階層逆上ったテキスト参照に関して、最初から1階層ではなく、N階層（Nは1より大きなある値）の範囲を参照範囲と決めておく方法もあるし、また、1階層の範囲で曖昧性解消ができなかった時に、参照範囲を1階層ずつ増やしていくという方法もある。例えば、リンク元及びリンク先の前後一階層分のテキストを参照して曖昧性を解消できなかった時に、リンク先のリンク、若しくはリンク元のリンクをたどり、共起語を検出していく方法である。

【0066】また、単純に階層を増やしていくのではなく、1階層目はそのテキスト全体を参照するが、2階層目はリンク元キー文字列（或いはリンク元キー文字列の前後一定範囲のテキスト）のみを参照範囲とするような方法も考えられる。図10の例で言えば、テキスト23に対して、1階層目のテキスト22はその全体を参照し、2階層目のテキスト21はリンク元キー文字列の「ベ이스ターズファン」の部分のみを参照するという方法である。

【0067】また、リンクを逆上るだけでなく、リンク先の方向も含めて階層を増やしていく方法も考えられる。例えば、図10の例において、テキスト23に対して、リンク元のテキスト22へ1階層分逆上り、今度はそのリンク先の方向にたどって、テキスト24を参照することも可能である。テキスト23に対して、2階層分の範囲として、テキスト21、テキスト22、テキスト24の3つを参照する方法、リンク元方向のみにたど

てテキスト21とテキスト22の2つを参照する方法、あるいは、兄弟関係のリンク参照を優先してテキスト22とテキスト24の2つを参照する方法などが考えられる。

【0068】次に、図11のフローチャートを用いて、本実施の形態の動作を説明する。まず、図11のステップS201にあるように、読み込み手段2によって、ハイパーテキストデータベース1からテキストを読み込んで、入力メモリ3に書き込む。読み込み手段2の実現方法によっては、ステップS201と以降のステップSとを交互に進めるような処理手順もとる得るが、ここでは前述の読み取り手段2の実現方法（b）を用いて、ある範囲内のテキストをまとめて読み込むものとし、その結果、読み込まれたテキストの件数はN件であったとする。

【0069】次に、N件のテキストの各々について、ステップS204以降の手順を実行する。図11のフローチャートでは、ステップS202でkの値を1にセットした上で、ステップS209でkの値を1ずつ増やしながら、ループ処理を実行している部分が、これに該当する。k番目のテキスト（k=1～N）に対する処理として、まず、ステップS204を実行する。ステップS204では、単一テキスト解析手段4によってテキストkを解析し、テキストk内に出現する固有表現の候補を検出して、中間メモリ5へ書き込む。ここで検出された固有表現候補の数をMkとし、個々の固有表現候補をc[k, j]（j=1～Mk）で表すものとする。すなわち、c[k, j]は、テキストkにおいて検出されたj番目の固有表現候補である。

【0070】次に、テキストkにおいて検出されたMk個の固有表現候補の各々について、ステップS207以降の手順を実行する。図11のフローチャートでは、ステップS205でjの値を1にセットした上で、ステップS208でjの値を1ずつ増やしながら、ループ処理を実行している部分が、これに該当する。j番目の固有表現候補c[k, j]に対する処理として、ステップS207を実行する。ステップS207では、複合テキスト解析手段6によって、テキストkに対するリンク元あるいはリンク先のテキストを参照して解析し、固有表現候補c[k, j]の確からしさを計算し、出力メモリ7へ書き込む。読み込み手段2の実現方法によっては、ステップS207の段階で、テキストkのリンク元テキストやリンク先テキストが入力メモリ3に含まれていないというケースもあり得る。その場合は、そのようなテキストkに関する固有表現候補c[k, j]については単一テキスト解析手段4の結果をそのまま出力することにしてもよいし、あるいは、ステップS207の段階で、読み込み手段2によってリンク元あるいはリンク先のテキストを改めて読み込むことにしてもよい。

【0071】テキストkに対するMk個の固有表現候補

の確からしさの計算が終了したら、ステップ S 206 を
 経て、次のテキスト (k+1) の処理へ進む。そして、
 N 件のテキストに対する処理が完了したら、ステップ S
 203 を経て、フローチャート全体の処理が終了する。

【0072】次に、本発明に係る第 2 の実施形態につい
 て添付図面を参照しながら詳細に説明する。図 12 に
 は、本発明の固有表現判別装置を利用した地図検索装置
 の構成を示すブロック図が示されている。

【0073】図 12 に示されるように固有表現判別装置
 を利用した地図情報検索装置は、ハイパーテキストデー
 タベース 50、固有表現判別装置 51、位置依存コンテン
 ツデータベース 52、地図データベース 53、データ
 ベース検索装置 54、表示装置 55、位置条件入力装置
 56 を備える。

【0074】ハイパーテキストデータベース 50 は、ハ
 イパーテキストが格納されている。例えば、インターネ
 ット上の WWW がこれに相当する。

【0075】固有表現判別装置 51 は、ハイパーテキス
 トデータベース 50 内のテキストから地名表記を検出す
 る。これまで説明した第 1 の実施の形態がこれに相当す
 る。ただし、本実施形態は、固有表現のうち地名と判定
 されたもののみを使用する。

【0076】位置依存コンテンツデータベース 52 は地
 名表記と、そのノードページ番号が対応付けられて格納
 されている。例えば、「東京都港区」という地名表記に
 関して、ノードページ 31 が対応し、「群馬県前橋市」
 という地名表記に関して、ノードページ 40 が対応して
 いる。

【0077】地図データベース 53 は、地図の 2 次元座
 標データと、その上にマッピングされた地名表記を格納
 している。

【0078】位置条件入力装置 56 は、「東京都港区」
 というような地名表記を利用者が入力するための装置で
 ある。キーボードのような文字列入力手段、マウスなど
 のポインティングデバイス、さらには、GPS のよう
 な人や車の現在位置を自動的に取得するシステムなどが
 用いられる。

【0079】データベース検索装置 54 は、位置条件入
 力装置 56 で指定された条件で、地図データベース 53
 と位置依存コンテンツデータベース 52 を検索して、そ
 の結果を表示装置 55 に表示する。

【0080】図 14 は、図 13 の位置依存コンテンツデ
 ータベースを用いて地図上にコンテンツを表示した例で
 ある。位置条件入力装置 56 からは関東地方の地名が指
 示されたものとする。

【0081】次に、本発明に係る第 3 の実施形態につい
 て図面を参照して詳細に説明する。図 15 を参照する
 と、本発明に係る第 3 の実施形態は、入力装置 100、
 データ処理装置 110、記憶装置 120、出力装置 14
 0 を備え、さらに、上述した第 1 の実施形態の固有表現

判別装置を実現するためのプログラムを記録した記録媒
 体 130 を備える。この記録媒体 130 は、磁気ディス
 ク、半導体メモリ、CD-ROM その他の記録媒体であ
 ってよい。

【0082】入力装置 100 は、マウス、キーボード
 等、操作者からの指示を入力するための装置である。ま
 た、出力装置 140 は、表示画面、プリンタ等のデータ
 処理装置 110 による処理結果を出力する装置である。

【0083】固有表現判別装置を実現するためのプログ
 ラムは、記録媒体 130 からデータ処理装置 110 に読
 み込まれ、データ処理装置 110 の動作を制御し、記憶
 装置 120 に入力メモリ 3 と中間メモリ 5 と出力メモリ
 7 とを生成する。データ処理装置 110 は、固有表現判
 別装置を実現するためのプログラムの制御により第 1 の
 実施形態における読み込み手段 2、複合テキスト解析手
 段 6、及び単一テキスト解析手段 4 による処理と同一の
 処理を実行する。

【0084】次に、本発明に係る第 4 の実施形態につい
 て図面を参照して詳細に説明する。図 16 を参照する
 と、本発明に係る第 4 の実施形態は、入力装置 200、
 データ処理装置 210、記憶装置 220、出力装置 24
 0 を備え、さらに、上述した第 2 の実施形態の地図情報
 検索装置を実現するためのプログラムを記録した記録媒
 体 230 を備える。この記録媒体 230 は、磁気ディス
 ク、半導体メモリ、CD-ROM その他の記録媒体であ
 ってよい。

【0085】入力装置 200 は、マウス、キーボード
 等、操作者からの指示を入力するための装置である。ま
 た、出力装置 240 は、表示装置、プリンタ等のデータ
 処理装置 210 による処理結果を出力する装置である。

【0086】地図情報検索装置を実現するためのプログ
 ラムは、記録媒体 230 からデータ処理装置 210 に読
 み込まれ、データ処理装置 210 の動作を制御し、記憶
 装置 220 に入力メモリ 3、中間メモリ 5、出力メモリ
 7、位置依存コンテンツデータベース 52、地図デー
 タベース 53 を生成する。データ処理装置 210 は、地図
 情報検索装置を実現するためのプログラムの制御により
 第 1 の実施形態における読み込み手段 2、複合テキスト
 解析手段 6、単一テキスト解析手段 4、第 2 の実施形態
 におけるデータベース検索装置 54 による処理と同一の
 処理を実行する。

【0087】上述した実施形態は、本発明の好適な実施
 の形態である。但し、これに限定されるものではなく、
 本発明の要旨を逸脱しない範囲内において種々変形実施
 が可能である。

【0088】

【発明の効果】以上の説明より明らかなように本発明
 は、固有表現の検出・判別に関して、その固有表現の出
 現したテキスト内だけではなく、リンク元やリンク先の
 テキストも参照して解釈することで、従来よりも曖昧性

を解消することが可能になり、高精度な固有表現判別を実現できる。

【図面の簡単な説明】

【図 1】本発明の実施の形態を示すブロック図である。

【図 2】本発明の実施の形態における処理対象の例を示す図である。

【図 3】HTML で記述されたテキストの例を示す図である。

【図 4】中間メモリの管理テーブルの例を示す図である。

【図 5】本発明の実施の形態における処理対象の例を示す図である。

【図 6】複合テキスト解析手段の構成を表すブロック図である。

【図 7】固有表現辞書の内容を示す図である。

【図 8】本発明の実施の形態における処理対象の例を示す図である。

【図 9】本発明の実施の形態における処理対象の例を示す図である。

【図 10】本発明の実施の形態における処理対象の例を示す図である。

【図 11】本発明の実施形態の動作を示すフローチャートである。

【図 12】本発明に係る第 2 の実施の形態を示すブロック図である。

10

* 【図 13】位置依存コンテンツデータベースのデータ内容を示す図である。

【図 14】地図情報検索装置の表示結果を示す図である。

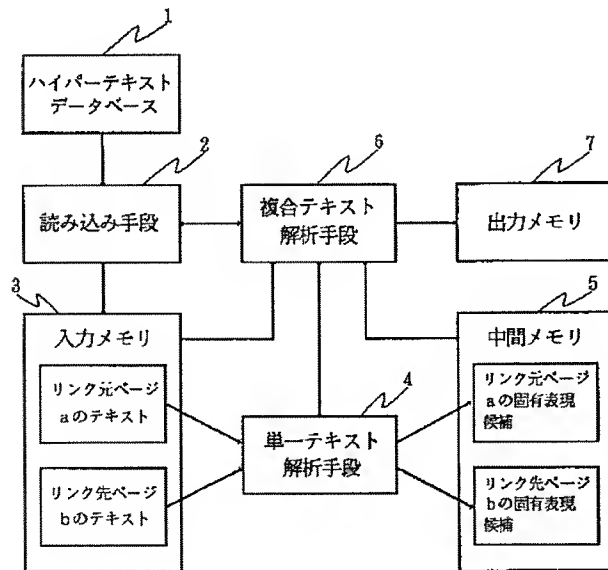
【図 15】本発明に係る第 3 の実施の形態を示すブロック図である。

【図 16】本発明に係る第 4 の実施の形態を示すブロック図である。

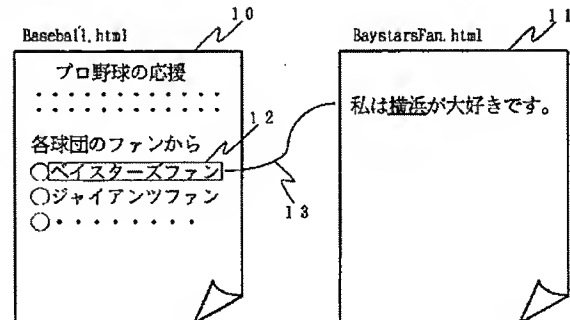
【符号の説明】

- 1 ハイパーテキストデータベース
- 2 読み込み手段
- 3 入力メモリ
- 4 単一テキスト解析手段
- 5 中間メモリ
- 6 複合テキスト解析手段
- 7 出力メモリ
- 10 リンク元ページ
- 11 リンク先ページ
- 12 リンク元キー文字列
- 13 リンク
- 14 リンク元キー文字列の前後一定範囲のテキスト
- 40 固有表現辞書における表記
- 41 固有表現辞書における種別
- 42 固有表現辞書における共起語

【図 1】



【図 2】



【図 4】

検出対象テキスト識別情報			
識別番号	固有表現候補	位置情報	種類

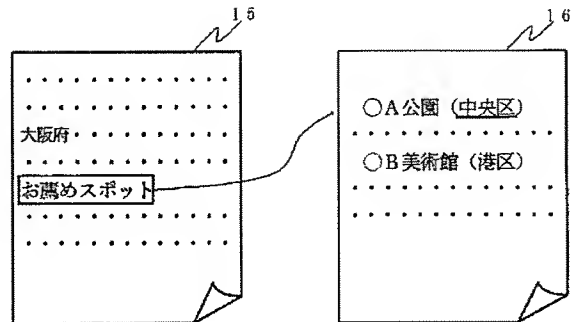
【図3】

Baseball.html

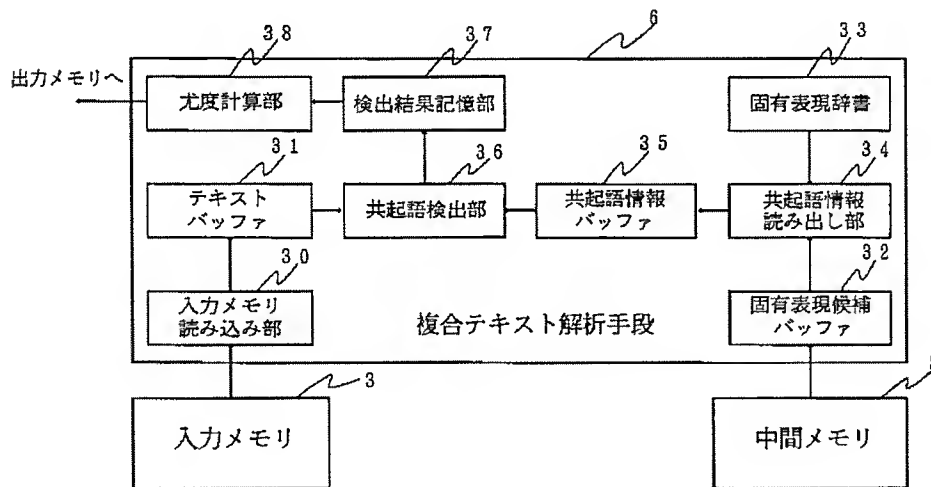
```

<HTML>
<HEAD>
<TITLE> Baseball Fan's Page </TITLE>
</HEAD>
<BODY>
<H1> プロ野球の応援 </H1>
.....
各球団のファンから
○<A HREF="BaystarsFan.html">ベイスターズファン</A>
○<A HREF="GiantsFan.html">ジャイアンツファン</A>
○.....
</BODY>
</HTML>
  
```

【図5】



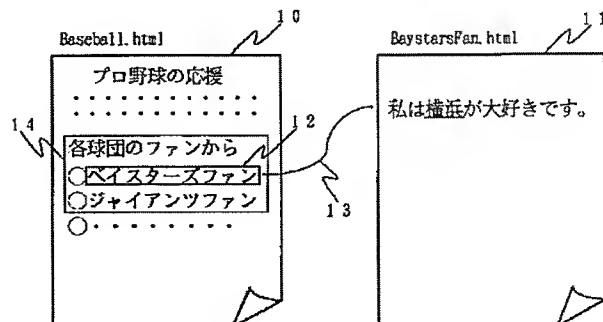
【図6】



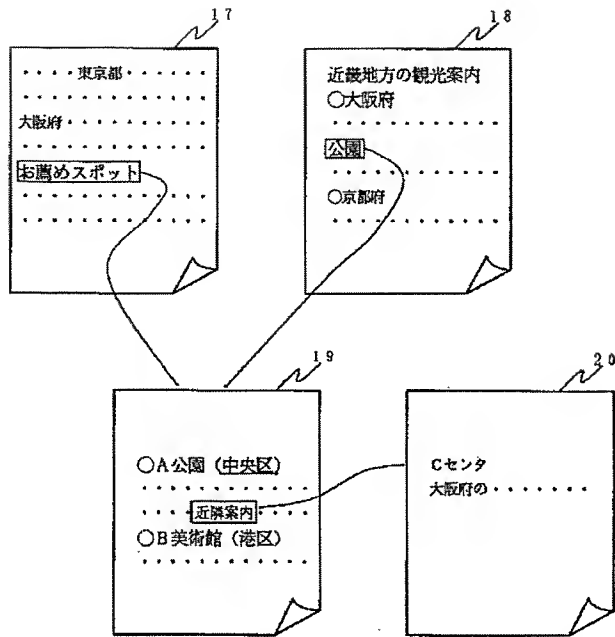
【図7】

固有表現	種別	共起語リスト
横浜	地名	神奈川、関東、.....
	組織名	ベイスターズ、プロ野球、球団.....
⋮	⋮	⋮

【図8】



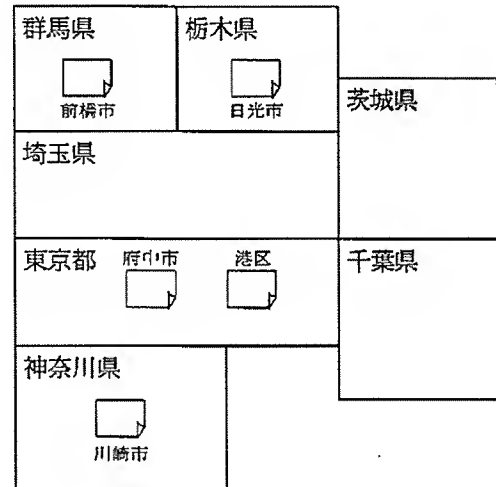
【図 9】



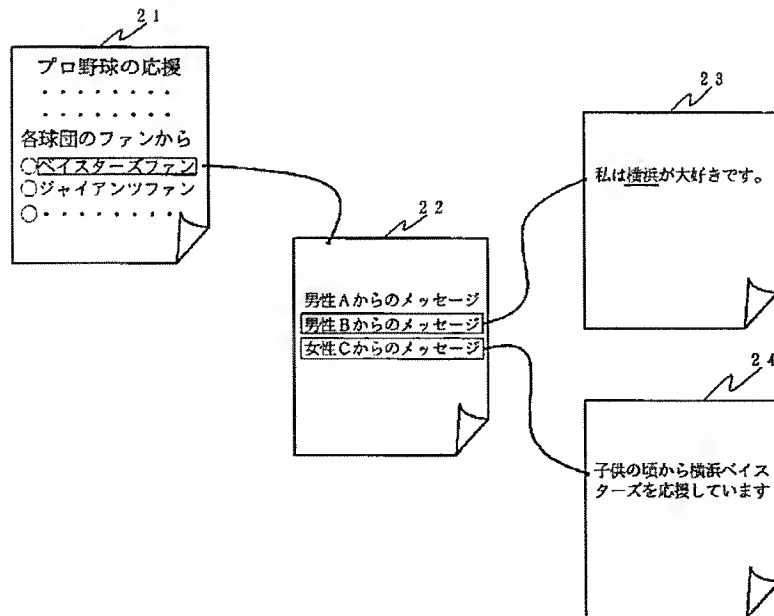
【図 13】

地名表記	ノードページ番号
東京都港区	31
東京都府中市	39
神奈川県川崎市	40
群馬県前橋市	40
栃木県日光市	52

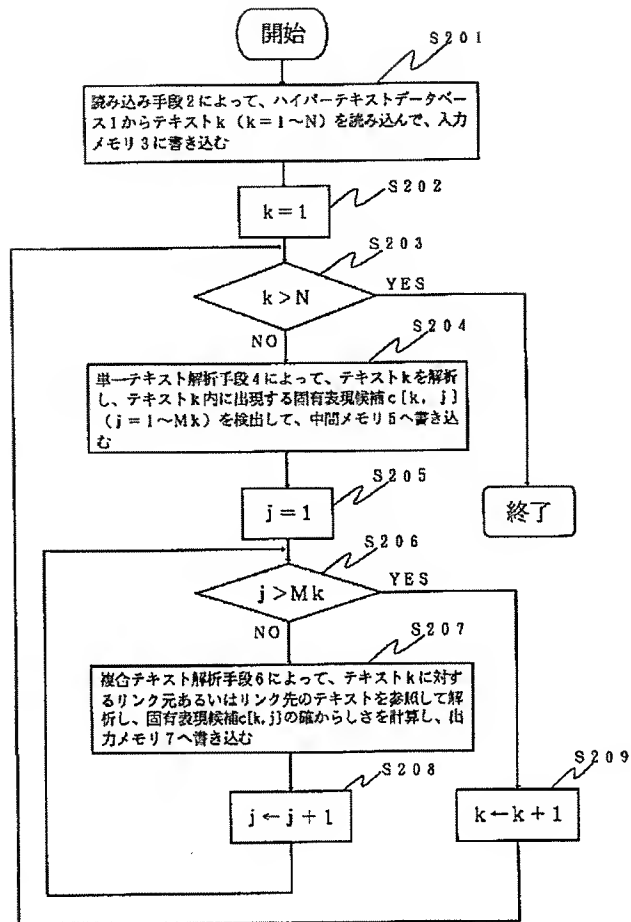
【図 14】



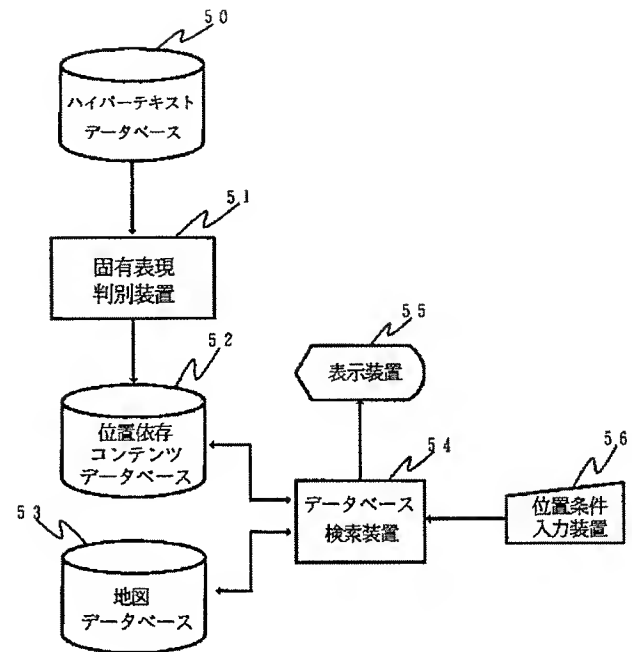
【図 10】



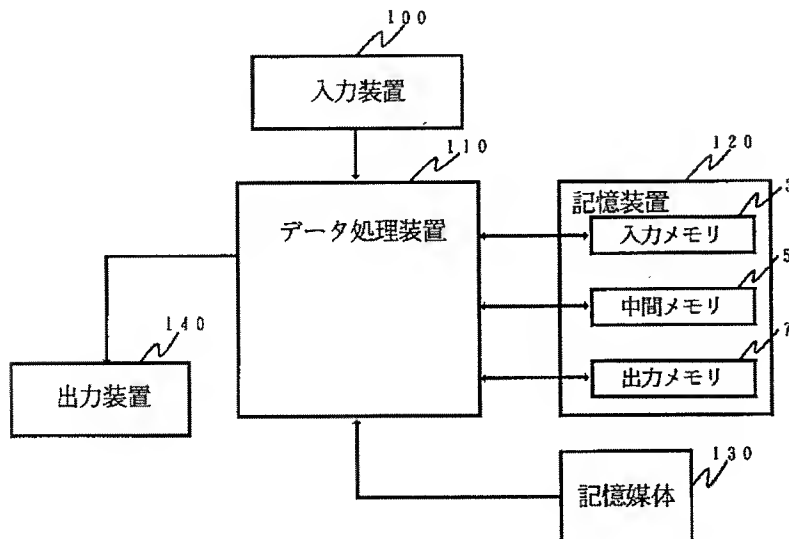
【図11】



【図12】



【図15】



【図 16】

